# MoE-SPNet: A Mixture-of-Experts Scene Parsing Network

Huan Fu[a], Mingming Gong[b], Chaohui Wang[c], Dacheng Tao[a]

[a]*UBTECH Sydney AI Centre, SIT, FEIT, The University of Sydney, J12 Cleveland St, Darlington NSW 2008, Australia*
[b]*Department of Biomedical Informatics University of Pittsburgh Cublicle 520c, 5607 Baum Bouevard, Pittsburgh, PA 15206, America*
[c]*Laboratoire d'Informatique Gaspard Monge - CNRS UMR 8049, Université Paris-Est, 77454 Marne-la-Vallée Cedex 2, France*

## Abstract

Scene parsing is an indispensable component in understanding the semantics within a scene. Traditional methods rely on handcrafted local features and probabilistic graphical models to incorporate local and global cues. Recently, methods based on fully convolutional neural networks have achieved new records on scene parsing. An important strategy common to these methods is the aggregation of hierarchical features yielded by a deep convolutional neural network. However, typical algorithms usually aggregate hierarchical convolutional features via concatenation or linear combination, which cannot sufficiently exploit the diversities of contextual information in multi-scale features and the spatial inhomogeneity of a scene. In this paper, we propose a mixture-of-experts scene parsing network (*MoE-SPNet*) that incorporates a convolutional mixture-of-experts layer to assess the importance of features from different levels and at different spatial locations. In addition, we propose a variant of mixture-of-experts called the adaptive hierarchical feature aggregation (*AHFA*) mechanism which can be incorporated into existing scene parsing networks that use skip-connections to fuse features layer-wisely. In the proposed networks, different levels of features at each spatial location are adaptively re-weighted according to the local structure and surrounding contextual information before aggregation. We demonstrate the effectiveness of the proposed methods on two scene parsing datasets including PASCAL VOC

---

*Email addresses:* hufu6371@uni.sydney.edu.au (Huan Fu), gongmingnju@gmail.com (Mingming Gong), chaohui.wang@u-pem.fr (Chaohui Wang), dacheng.tao@sydney.edu.au (Dacheng Tao)

2012 and SceneParse150 based on two kinds of baseline models FCN-8s and DeepLab-ASPP.

## 1. Introduction

Scene parsing or semantic image segmentation, which predicts a category-level label (such as "sky", "dog" or "person") for each pixel in a scene, is an important component in scene understanding. A perfect parsing can contribute to a variety of applications including unmanned vehicles, environmental reconstruction, and visual SLAM. Many other fundamental computer vision problems can benefit from the parsing of an image, such as medical image analysis, tracking, and object detection [1, 2, 3]. However, scene parsing is a very challenging high-level visual perception problem as it aims to simultaneously perform detection, reconstruction, segmentation, and multi-label categorizing [4, 5].

Since feature representation is critical to pixel-level labeling problems, classical methods focus on designing handcrafted features for scene parsing [6]. Since the hand-crafted features alone can only capture local information, probabilistic graphic models such as conditional random fields (CRFs) are often built on these features to incorporate smoothness or contextual relationships between object classes [7]. Recently, deep learning approaches such as deep convolutional neural networks (DCNNs) have earned immense success in scene parsing. In particular, fully convolutional networks (FCNs)-based approaches have demonstrated promising performance on several public benchmarks [5, 8, 9, 10, 11, 12].

A common strategy adopted in all the CNN-based methods is to aggregate multi-scale/level features from multiple CNN layers [5] or from a specific layer [8], which is a key component to obtain high-quality dense predictions because the multi-level features capture different levels of abstractions of a scene. The standard way to combine hierarchical features/predictions is to either concatenate multi-level features [13, 14, 15, 16, 17, 18, 19, 20] or equivalently aggregate the prediction maps by average

2

pooling [5]. However, the linear feature aggregation methods are not able to evaluate the relative importance of the semantic and spatial information in each level of features. The information at different scales is complementary because the higher-level convolutional features contain larger-scale contextual information which is beneficial for classification, while the lower-level features have higher spatial resolution which produces finer segmentation masks [21]. The information at different scales is also complementary since they are from different receptive fields. There is thus a trade-off between the semantic and the spatial information. In addition, the average pooling ignores the spatial inhomogeneity of a scene, which is improper since different objects may prefer features from different scales/levels. For example, textured objects such as "grass" and "trees" can be easily distinguished from lower-level features while texture-less objects like "bed" and "table" require higher-level features to capture the global shape information.

In this paper, we propose a mixture-of-experts [22] scene parsing network (*MoE-SPNet*) which learns to aggregate multi-level convolutional features according to the image structures. Specifically, we treat each network branch that contains a specific level/scale of features/predictions as an expert and aggregate them using the weights generated by a trainable convolutional gating network. The gating network also has a convolutional architecture and outputs a weight map for the entire image. The proposed MoE-SPNet is motivated by the following three observations: 1) The lower-level convolutional features contain more precise boundary information but tend to yield more incorrect predictions, while the higher-level features contain more contextual and semantic information but less spatial information. 2) Different levels/scales of features reflect the visual properties of different-sized objects because they are extracted by receptive fields with different sizes. Notably, small objects are more likely to be misclassified to their background if using higher-level features because larger receptive fields introduces much noise to these small objects. 3) The relative importance of different levels of features varies with spatial location; it relies on the local image structure and surrounding contextual information. Obviously, a linear combination of these features by average pooling cannot capture the homogeneity of a scene and assess the importance of different feature levels. On the contrary, the proposed MoE-SPNet

overcomes the limits of linear combination by aggregating different level of features in a nonlinear and adaptive way.

Since MoE-SPNet is only able to adaptively aggregate multi-scale features generated from a single CNN layer, we further propose a variant of MoE called adaptive hierarchical feature aggregation scheme (AHFA) which can be incorporated into the existing parsing networks that aggregate hierarchical features using skip-connections. For example, the original FCN architecture combines features from the last convolutional layer with previous layers by successive upsampling and aggregation. Employing AHFA will enable the parsing networks such as FCN to learn weights at each stage and aggregate the features adaptively as done in MoE-SPNet. In this paper, we focus on exploiting AHFA for the original FCN, leading to a new network architecture denoted as *FCN-AHFA*.

We demonstrate the effectiveness of our MoE-SPNet and FCN-AHFA on two challenging benchmarks for scene parsing, PASCAL VOC 2012 [23] and SceneParse150 [24], and achieve the state-of-the-art or comparable results. Also, the experimental results show that our MoE-SPNet and FCN-AFHA consistently improve the performance of all the evaluated baseline networks, and thus demonstrate the value of the proposed methods. In addition, the produced weight maps can help us understand the reason that some image structures prefer higher-level convolutional features while others prefer lower-level features.


## 2. Related work

Segmentation is a fundamental problem in scene understanding. While some works focus on low-level segmentation which segments a scene into some regions that share certain characteristics or computed property, such as color, intensity, or texture [25, 26, 27, 28], high-level segmentation (scene parsing or semantic segmentation), which assigns a category-level label to each pixel of a scene, receives much attention recently.

In the past decade, the successful scene parsing methods rely on handcrafted local features like colour histogram and textons [6, 29, 30, 31, 32, 33], and shallow classifiers such as Boosting [6, 34], Random Forests [35, 36], Support Vector Machines [37]. Due

to the limited discriminative power of local features, a lot of efforts have been put into developing probabilistic graphical models such as CRFs to enforce spatial consistency and incorporate rich contextual information [38, 7, 39, 40]. Recently, deep learning methods typified by DCNNs have achieved state-of-the-art performance on various computer vision tasks, such as image classification and multi-class object detection.

Also, the DCNN architectures such as VGG [41] and ResNet [42] originally developed for image classification have been successfully transferred to scene parsing. Specifically, Long *et al.* [5] proposed the fully convolutional network (FCN) which applied DCNNs to the whole image and directly produced dense predictions from convolutional features, making it possible to get rid of bottom-up segmentation steps [43] and train the parsing network in an end-to-end fashion.

The impressive performance of FCNs is largely due to the aggregation of multi-level or multi-scale features/predictions. There are mainly two types of aggregation methods: share-nets and skip-nets [44]. The skip-nets, which merge multi-level features/predictions from a single network, are computationally more efficient than the share-nets. Furthermore, they have been refined to enable end-to-end training by normalizing the features from different levels. For example, Hariharan *et al.* [4] concatenated the multi-level features together after certain normalization methods like L2 normalization. However, the concatenation of hierarchical features results in high-dimensional features and is thus time-consuming. The FCN-8s [5] model aggregated features from the last three convolutional blocks by averagely pooling over layers. Similarly, Chen *et al.* [45] combined the features which were extracted by applying multi-layer perceptrons on the original image and the pooling layers. However, linear combination of multi-scale features does not sufficiently exploit the geometric properties, contextual information, and the spatial-semantic tradeoff. Recently, Ghiasi *et al.* [21] found that directly summing up multi-scale features cannot achieve desirable results, as the learned parameters tended to down-weight the contribution of lower-level features (higher resolution) to suppress the effects of noisy predictions. They proposed the laplacian pyramid refinement approach which computed a boundary mask from higher-level semantic predictions to filter out the noisy predictions in lower-level features. However, we aim to learn the mask weights from multi-level features instead of

calculating a boundary mask by manually designed mathematical operations.

Share-nets combine features from shared networks built on multiple rescaled images. For example, Farabet *et al.* [43] transformed the raw image through a laplacian pyramid, and each level of which was fed into a CNN. The produced sets of feature maps of all scales were concatenated to form the final representation. Similarly, Lin *et al.* [46] resized the original image to three scales and concatenated the multi-scale features. Aside from concatenation, average pooling [47] and max pooling [48] were adopted over scales to merge multi-scale features. However, average or max pooling either treats the multi-scale features equally or losses too much information. Targeting this problem, Chen *et al.* [44] proposed the scale attention method which uses the attention model [49] over scales to focus on the features from the most relevant scales. Instead of aggregating multi-scale features at one time, Pinheiro *et al.* [50] proposed a multi-stage approach which fed multi-scale images successively to a recurrent convolutional neural network. Although the share-nets obtain much better performance, they are computationally more expensive than the single scale networks. Most recently, Chen *et al.* [8] developed an atros spatial pyramid pooling strategy (a variant of the share-nets) which extracted multi-scale features in a single network. However, the multi-scale features were still aggregated via an average pooling, and the performance had some gaps against the typical share-nets.

In this paper, we investigate how to adaptively aggregate multi-level or multi-scale features in a single network to further improve their performance and obtain deeper understanding of the special properties of the features from different layers. Specifically, we treat the network branches which obtain multi-level/scale features as expert networks, and propose MoE-SPNet, which learns some pixel-wise gating weight maps for each experts, to adaptively aggregating these features for a better solution for scene parsing. We also propose the AHFA scheme to further improve existing skip-nets [5, 21] that use stage-wise aggregation of hierarchical features. Since most of current parsing networks follow the similar forms as FCN or DeepLab, we can conclude that our technique is widely applicable.

6

## 3. Background

In this section, we first review the mixture-of-experts (MoE) framework and then review two typical scene parsing networks that employ fully convolutional architectures, i.e., FCN-8s [5] and DeepLab-ASPP [8].
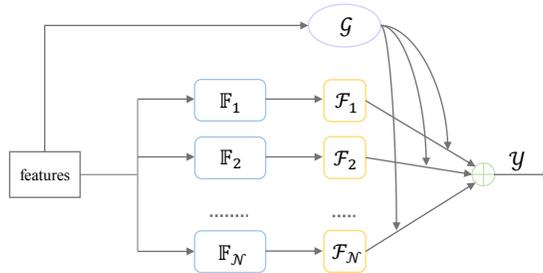
### 3.1. Mixture-of-Experts



Figure 1: **Mixture-of-Experts.** An illustration of mixture-of-experts. The same input is fed to different experts, resulting in different solutions for the whole problem space. Typically, the gating network $\mathcal{G}$ also receives the same input as the experts, and the weights are often nomalized by the *softmax* function. Here, $\mathbb{F}_i$ and $\mathcal{F}_i$ are learned intermediate features and a prediction correspond to expert $i$ respectively.

Mixture-of-experts [22] is one of the effective machine learning techniques which aims to adaptively aggregating multiple decisions from different experts. As shown in Fig. 1, MoE contains two key components: multiple correlated experts and a gating network. The multiple correlated experts are expected to learn the distribution specialized on a stochastic subspace of the whole problem space, and are thus complementary to each other. The gating network aims at learning weights for each expert according to their local efficiency. It should be noted that, the weights in the gating network are dynamically determined by the input features. Here, we take mixture-of-expert networks as an example, and introduce the conventional MoE with respect to two different error functions in the learning process.

### 3.1.1. Cooperation Encouraged Error Function

The error function which encourages cooperation among local experts exhibits the following form:

$$E_{coop} = \|y - \sum_{i=1}^{N} g_i o_i\|^2 \tag{1}$$

where $y$ is the target vector, $N$ is the number of experts, $o_i$ is the output of expert $i$, and $g_i$ from the gating network ($g_1 + g_2 + ... + g_N = 1$) represents the contribution of expert $i$ for the final prediction.

With this error function, the blend of the outputs from each expert is directly compared with the target, meaning that the parameters in each expert are updated according to the overall ensemble error. The strong coupling in the learning process encourages all the experts cooperating nicely, but tends to make each expert generalize to the whole problem space rather than to different subspaces of the whole problem space. Thus, the learned model via this error function may become inconsistent with the localization of the experts.

### 3.1.2. Competition Encouraged Error Function

Addressing the shortage in cooperation encouraged error function, Jacobs *et al*. [51] defined a competition encouraged error function as:

$$E_{comp} = \sum_{i=1}^{N} g_i \|y - o_i\|^2 \tag{2}$$

From the definition, this error function actually measures the expected value of differences between the target and each local experts. Thus, each expert directly responds to their own occasions and obtain a complete output over the whole problem space instead of a residual. After the training process, a single expert prefer to generate a solution for a specific training case, and the gating network here plays a role in selecting one or several experts for a given input. In this case, the experts still have some indirect coupling of each other due to the gating network.

### 3.2. FCN and DeepLab-ASPP

FCN-8s [5] applies a deep convolutional architecture, *e.g.* VGG net [41], in a fully convolutional fashion to extract hierarchical features with different strides (32x, 16x,

and 8x), and combine these features stage by stage from a deeper (coarser) layer to a shallower (finer) layer. Specifically, built on the 16-layer VGG (VGG16) architecture, FCN-8s replaces fully-connected layers with convolutional layers to generate the prediction feature maps with stride 32.

DeepLab-ASPP reduces the stride of 32x feature maps of FCN to 8 by using dilated convolutions (atrous algorithm) [52], which introduces zeros to increase the convolution fields for the convolutional kernels. Then, the atrous spatial pyramid pooling (ASPP) strategy, which employs multiple parallel filters with different dilation rates on the `pool5` layer, is adopted to exploit multi-scale features. The generated predictions from the multi-scale features are simply summed together to produce the final prediction.

## 4. Approach

In this section, we first present how to incorporate MoE in a scene parsing network and describe the details of the proposed MoE-SPNet. Second, we introduce adaptive hierarchical feature aggregation (AHFA) scheme which is a variant of MoE and show how it can incorporated into the skip-net FCN-8s [5] to form a new network FCN-AHFA. The AHFA scheme can be incorporated into other skip-nets in a similar way.

### 4.1. MoE-SPNet

We develop a mixture-of-experts scene parsing network (MoE-SPNet) which aims to learn predictions by considering features computed with different receptive fields (experts) and adaptively aggregate these predictions (gating network) to produce final semantic segmentation masks. Our network is built on DeepLab-ASPP [8] which exploits different receptive fields for scene parsing.

Each expert in MoE-SPNet targets at learning a parsing mask from a specific receptive field. In particular, an expert adopts a dilated convolutional layer with a specific dilation rate ($e_i^1$) to obtain local structural and contextual information from features computed with a specific receptive field on top of the `pool5` layer. Followed by two additional convolutional layers with the filter size of $1 \times 1$ ($e_i^2$ and $e_i^3$), each expert can
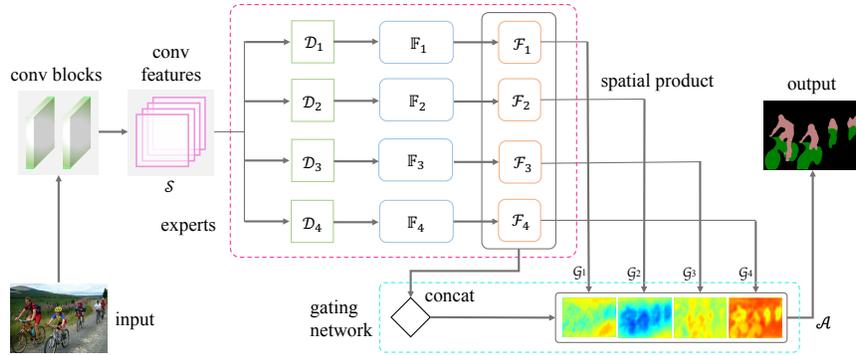
Figure 2: **MoE-SPNet.** An illustration of the proposed MoE-SPNet. $\mathcal{D}_i$ represents a dilated convolutional layer with a specific dilation rate. We learn 4 experts in this paper with the dilation rates of 6, 12, 18, and 24 respectively. Each expert learns a richer representation ( $\mathbb{F}_i$) of the input scene, and produces a solution (denote as $\mathcal{F}_i$) for the parsing task. $\mathcal{G}_i$ represents the weight map produced by the gating network for each parsing solution $\mathcal{F}_i$. $\mathcal{A}$ is the final segmentation mask which is the weighted aggregation of all $\mathcal{F}_i$.
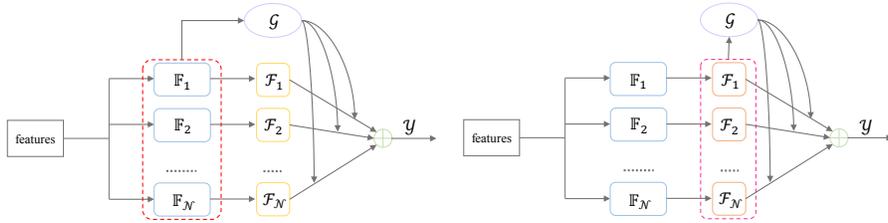


Figure 3: **Variants of MoE.** Left: Learning the gating network from high-level features $\mathbb{F}_i$ of each expert. Right: Learning the gating network from the predictions $\mathcal{F}_i$ of each expert.

learn a richer representation (denote as $\mathbb{F}_i$) of the input scene, and produce a solution (denote as $\mathcal{F}_i$) for the parsing task. Thus, with different dilation rates, the network can obtain some experts corresponding to different parsing solutions. Specifically, each experts are supervised by the ground-truth parsing via *softmax regression*. Thus, each channel of $\mathcal{F}_i$ corresponds to the probability of belonging to a category.

As shown in Fig. 2, the gating network in our MoE-SPNet is different from the standard gating network that learns weights from the input features to combine a series
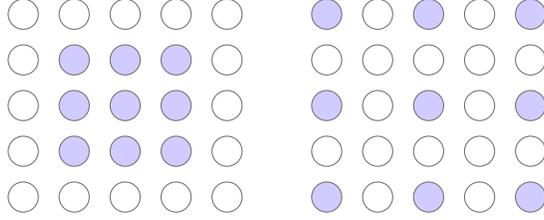
Figure 4: **Dilation.** Convolutional layers with the kernel size of 3. Left: standard convolutional layer. Right: dilated convolutional layer with the dilation rate of 2.

of classifiers. We learn the gating network from the segmentation maps generated by different experts instead of the same input features fed to the experts. Fig. 3 shows two variants of the standard MoE: one uses the high-level features $\mathbb{F}_i$ and the other uses the predictions $\mathcal{F}_i$ for the gating network. These two variants are supposed to perform better than the standard MoE, because of adoption of higher-level representations for the gating network. Since prediction maps are directly supervised using ground-truth segmentation maps, they contain the richest semantic information and are most suitable for training the gating network. Another advantage of using predictions $\mathbf{F}_i$ to train the gating network is that the number of gating network parameters can be reduced, leading to lower computational and memory cost.

To train the gating network, we concatenate $\mathcal{F}_1$ to $\mathcal{F}_\mathcal{N}$ ($\mathcal{N}$ is the number of experts), denoted as $\mathcal{F}$, and learn a non-linear function via two convolutional layers ($g_1'$ and $g_2'$) from these features to the corresponding gating features, denoted as $\mathcal{G}$, which follows the form:

$$\mathcal{G} = (\mathcal{F} * \mathcal{K}_{g_1'}) * \mathcal{K}_{g_2'}, \tag{3}$$

where $\mathcal{K}_{g_1'}$ and $\mathcal{K}_{g_2'}$ are kernels of the convolutional layers $g_1'$ and $g_2'$ with the kernel size of $3 \times 3$ and $1 \times 1$ respectively, "$*$" is the convolution operator, and the gating features set $\mathcal{G}$ in this paper consists of $\mathcal{G}_1$ to $\mathcal{G}_\mathcal{N}$. Followed by a normalisation process, the weight located at $(i, j)$ for expert $l$ can be calculated as:

$$w_l(i, j) = \frac{e^{\mathcal{G}_l(i,j)}}{\sum_{k=1}^{\mathcal{N}} e^{\mathcal{G}_k(i,j)}}. \tag{4}$$

11

After obtaining the weight maps in the gating networks, each channel $\mathcal{F}_i$ is multiplied by the corresponding weight map $\mathcal{W}_i$, and the aggregated output can be calculated as:

$$\mathcal{A} = \sum_{i=1}^{\mathcal{N}} \mathcal{F}_i \otimes \mathcal{W}_i, \tag{5}$$

where "$\otimes$" denotes element-wise product in each channel.

We train MoE-SPNet using the cost function consisting of a cooperation encouraged error term and a weakened competition encouraged error term:

$$\mathcal{L} = \Phi(\mathcal{Y}, \mathcal{A}) + \sum_{i=1}^{N} \Phi(\mathcal{Y}, \mathcal{F}_i \otimes \mathcal{W}_i), \tag{6}$$

where the "$\Phi$" represents the multinomial logistic regression error. Note that, all of the experts in our parsing network are addressing the single occasion rather than different occasions, thus the competition between these experts should not be strong. As a result, we ignore the gating factors in the typical competition encouraged error term.

### 4.2. FCN-AHFA



Figure 5: **AHFA:** An illustration of the proposed adaptively hierarchical features aggregation(AHFA) technique, which is another variant of mixture of experts.

We investigate how to incorporate the mechanism of MoE into another popular parsing network architecture with stage-wise fusions of features from different layers. We hypothesize that the gating map for each expert can be directly learned from the expert itself and propose an adaptive hierarchical feature aggregation (AHFA) mechanism which is a variant of the proposed MoE in Sec. 4.1 by assuming sparse connections in

the gating network. We take the typical parsing network FCN-8s [5] as an example to demonstrate the effectiveness of AHFA.

To take advantage of contextual information, FCN-8s produces a finer 16x-prediction with 16 pixel stride (16x) by adding a $1 \times 1$ convolutional layer on top of the `pool4` layer. The 32x-prediction is then upsampled to the same size of the 16x-prediction via a learnable deconvolutional layer, and then summed up with the 16x-prediction to accomplish one stage of combination. Finally, the above combined prediction is further aggregated with higher resolution (8x) features by applying the same strategy. The final prediction with stride 8 is upsampled back to the input image resolution.



Figure 6: **AHFA for skip-nets.** The feature maps are fused by stage-wise combination in skip-nets. In each combination stage, we learn a soft weight map for each level of features followed by a weighted pooling step over adjacent levels.

Now we describe the details of AHFA in FCN-8s to adaptively merge hierarchical features (32x, 16x, and 8x), resulting in a modified model which we call FCN-AHFA. An illustration of FCN-AHFA is shown in Fig. 6. In the first stage, on top of the 32x-prediction, denoted as $\mathcal{F}^{32\text{x}} \in \mathbb{R}^{H \times W \times C}$, we add a convolutional layer with the kernel size of $3 \times 3$ and the stride of $1$ followed by a sigmoid layer to produce a dense probabilistic weight map $\mathcal{W}^{32\text{x}} \in \mathbb{R}^{H \times W}$. Here, $H$, $W$, and $C$ denote the height, width, and the number of channels of the 32x feature maps, respectively. Then the weight located at $(i, j)$ in $\mathcal{W}^{32\text{x}}$ can be calculated as:

$$w^{32\text{x}}_{(i,j)} = \frac{1}{1 + e^{-\sum_{c=1}^{C}(f^{32\text{x}}_c * k^{32\text{x}}_c)(i,j)}}, \tag{7}$$

where $f_c^{32x}$ represents the $c$-th channel of $\mathcal{F}^{32x}$, $k_c^{32x}$ is the corresponding convolutional kernel, and "$*$" is the convolution operator. The weight function in Eq. (7) can be made more complex by introducing more convolutional and activation layers. However, we have observed from the experimental results that learning more complex weight functions only slightly improves the performance. After obtaining the weight map, each channel of $\mathcal{F}^{32x}$ is multiplied by $\mathcal{W}^{32x}$, resulting in the weighted features $\mathcal{H}^{32x} \in \mathbb{R}^{H \times W \times C}$ of which each channel is:

$$h_c^{32x} = \mathcal{W}^{32x} \otimes f_c^{32x}, \tag{8}$$

where $\otimes$ represents Hadamard product or entrywise product. Likewise, we reweight the 16x-prediction $\mathcal{F}^{16x}$ by the learned weight $\mathcal{W}^{16x}$ to obtain $\mathcal{H}^{16x} \in \mathbb{R}^{2H \times 2W \times C}$. At the final step of this stage, $\mathcal{H}^{32x}$ is upsampled to have the same size of the $\mathcal{H}^{16x}$ and linearly combined with it to produce the 16x aggregated feature:

$$\mathcal{A}^{16x} = \mathcal{H}^{16x} \oplus (\mathcal{H}^{32x})^{\uparrow}, \tag{9}$$

where $(\bullet)^{\uparrow}$ is a 2x upsampling operation via bilinear interpolation and $\oplus$ denotes the summing operation in each spatial location.

The aggregation strategy for the second stage is similar to that used in the first stage but is applied on $\mathcal{A}^{16x}$ and $\mathcal{F}^{8x}$. Hence, the the $c$-th channel of $\mathcal{A}^{8x}$ can be calculated as:

$$a_c^{8x} = (\mathcal{W}^{8x} \otimes f_c^{8x}) \oplus (\mathbb{W}^{16x} \otimes a_c^{16x})^{\uparrow}, \tag{10}$$

where $\mathcal{W}^{8x}$ and $\mathbb{W}^{16x}$ are the learned probabilistic weight maps for $\mathcal{F}^{8x}$ and $\mathcal{A}^{16x}$, repectively.

**Remark** The fixed-size filters ($3 \times 3$) used for learning the weight maps are actually adaptive to the size of semantic areas in the input image, because the higher-layer feature maps have smaller size. For example, the spatial areas corresponding to the original image considered by $k_c^{32x}$ are four times larger than that considered by $k_c^{16x}$. Also, the weight map of a layer is learned only from the feature maps in that layer.

This is different from existing mixture-of-experts [22] or the attention models [44] which usually learn the weights from the concatenation of features maps from all layers. Our method simplifies the weight learning network based on the observation that the feature maps in one layer already contain rich information about the corresponding weight map. Finally, with the learned weight maps, different levels of features can be aggregated adaptively by considering the relative spatial-semantic tradeoff at each spatial location.

## 5. Experiments

To demonstrate the effectiveness of the proposed MoE-SPNet and FCN-AHFA methods, we compare our methods with the existing methods on two challenging datasets, *i.e.* PASCAL VOC 2012 [23] and SceneParse150 [24]. We first describe the experimental settings including evaluation protocols and detailed implementations, and then report the experimental results with discussions.

### 5.1. Experimental Setting

**Evaluation Metrics** Four common metrics for scene parsing are used in our experiments, *i.e.* pixel accuracy, mean accuracy, mean IoU, and weighted IoU. Pixel accuracy indicates the proportion of correctly classified pixels. Mean accuracy indicates the average of the proportion of correctly classified pixels for all classes. IoU indicates the average intersection-over-union between the predicted and ground-truth pixels over all classes. Weighted IoU indicates the IoU weighted by total pixel ratio of each class. Let $L$ be the number of classes of interest, $l_{ij}$ represents the number of pixels belonging to class $i$ predicted as class $j$, and $N = \sum_{i=1}^{L} \sum_{j=1}^{L} l_{ij}$ is the number of pixels. The four metrics are computed as follows:

- Pixel Acc. : $\frac{1}{N} \sum_{i=1}^{L} l_{ii}$

- Mean Acc. : $\frac{1}{L} \sum_{i=1}^{L} \frac{l_{ii}}{\sum_{j=1}^{L} l_{ij}}$

- Mean IoU : $\frac{1}{L} \sum_{i=1}^{L} \frac{l_{ii}}{-l_{ii} + \sum_{j=1}^{L} (l_{ij} + l_{ji})}$

- Weighted IoU : $\frac{1}{N} \sum_{i=1}^{L} \frac{l_{ii} \sum_{j=1}^{L} l_{ij}}{-l_{ii} + \sum_{j=1}^{L} (l_{ij} + l_{ji})}$

15

It should be noted that pixel accuracy is biased to reflect the "stuff" categories such as grass and sky as they occupy more pixels. Instead, IoU is a more accurate measure of the classification performance on "things" categories such as person and car.

**Implementation** Since the proposed methods rely on semantic predictions in each level, our framework are trained in two stages. In the first stage, we train the basic network without MoE to produce hierarchical features containing semantic information. In the second stage, we add the gating network of MoE to the pre-trained baseline network and fine-tune the whole parsing network in an end-to-end fashion. For fair comparison, we also fine-tune the baseline network with the same iterations. We initialise the base convolutional architecture via the pre-trained VGG16, ResNet-50, and ResNet-101 [42] classification models on ILSVRC [53]. The fine-tuning stage follows a polynomial decay with the power of 0.9, the momentum of 0.9, and the weight decay of 0.0005. We implement our networks based on *Caffe* [54], and train them using 4 TITAN X GPUs with 12GB of memory per GPU. The batch size is set to 8 in all the experiments.

**Data Augmentation** Data augmentation techniques are used when training the parsing networks, which can be summarised as follows: 1). The training images are resized by the scaling factors: 0.5, 0.75, 1.0, 1.25, 1.5. 2). We randomly flip the training images horizontally. 3). The input samples of the models are randomly cropped from the training images with a fixed size.

*5.2. Benchmark Performance*

*5.2.1. PASCAL VOC 2012*

PASCAL VOC 2012 [55], which consists of 20 common object categories and one background category, is a well-known benchmark for semantic segmentation. The images contained in this dataset are split into three parts, including 1464 training images, 1449 validation images, and 1456 test images. Following [56], the training data with ground-truth segmentation masks are augmented to 10,582 images using the extra annotated images for VOC 2012. Since PASCAL VOC 2012 is an object-

| Image | GT | FCN-8s [5] | FCN-AHFA | Deeplab-ASPP [8] | MoE-SPNet |

Figure 7: **PASCAL VOC 2012 results.** A comparison of proposed MoE based parsing networks, *i.e.* FCN-AHFA, MoE-SPNet, with their baseline models, *i.e.* FCN-8s, DeepLab-ASPP. (Best view in colour.)

level segmentation benchmark, and each image in this dataset follows a simple foreground/background form. We only adopt mean IoU, which is a stricter and more convincing metric for scene parsing, to evaluate different methods following previous works.

In Tab. 1, we report our scores on the test server in different conditions to make a comparison with previous works. Our MOE-SPNet achieves about 2.5% improvement on the baseline model Deeplab-ASPP based on ResNet, and obtains comprisable results with current state-of-the-art algorithms on different settings. Also, our FCN-AHFA significantly outperforms the most typical baseline model FCN-8s, nearly closing the gap to current state-of-the-art methods. Fig. 7 gives qualitative comparison of different methods on several images.

|       |       |       |       |
| ----- | ----- | ----- | ----- |
| Image | GT    | FCN-8s [5] | FCN-AHFA |

|       |       |       |       |
| ----- | ----- | ----- | ----- |
| Image | Ground Truth | Deeplab-ASPP [8] | MoE-SPNet |

Figure 8: **SceneParse150 results.** The top part shows the segmentation results of FCN-8s [5] without or with our AHFA technique. The bottom part shows the 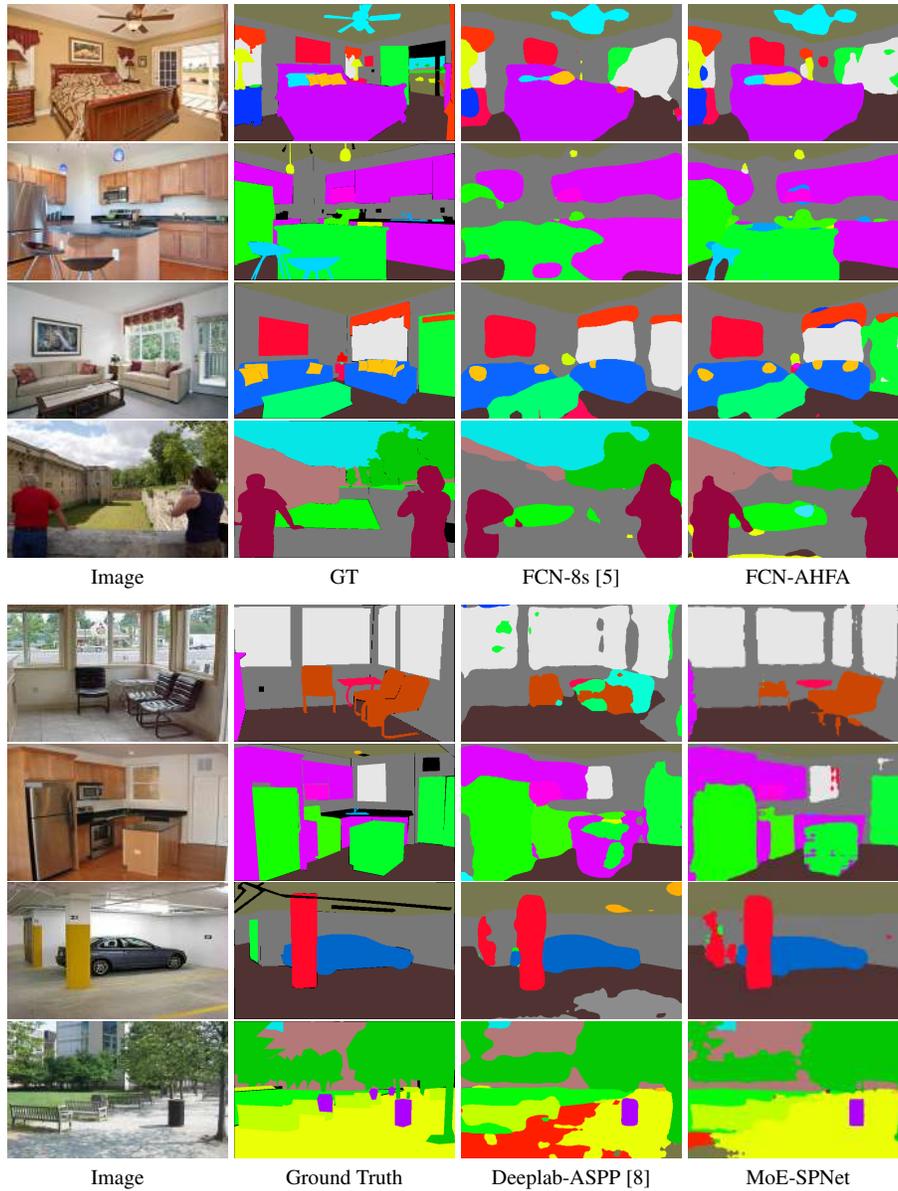segmentation results of Deeplab-ASPP [8] (with atrous spatial pyramid pooling) and our MOE-SPNet. (Best viewed in colour)

| | **mean** | areo | bike | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | mbike | person | plant | sheep | sofa | train | tv |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| VGG + PASCAL VOC | | | | | | | | | | | | | | | | | | | | | |
| SegNet [19] | 59.9 | 73.6 | 37.6 | 62.0 | 46.8 | 58.6 | 79.1 | 70.1 | 65.4 | 23.6 | 60.4 | 45.6 | 61.8 | 63.5 | 75.3 | 74.9 | 42.6 | 63.7 | 42.5 | 67.8 | 52.7 |
| FCN-8s [5] | 62.2 | 76.8 | 34.2 | 68.9 | 49.4 | 60.3 | 75.3 | 74.7 | 77.6 | 21.4 | 62.5 | 46.8 | 71.8 | 63.9 | 76.5 | 73.9 | 45.2 | 72.4 | 37.4 | 70.9 | 55.1 |
| Hypercolumn [13] | 62.6 | 68.7 | 33.5 | 69.8 | 51.3 | 70.2 | 81.1 | 71.9 | 74.9 | 23.9 | 60.6 | 46.9 | 72.1 | 68.3 | 74.5 | 72.9 | 52.6 | 64.4 | 45.4 | 64.9 | 57.4 |
| Zoom-out [57] | 69.6 | 85.6 | 37.3 | 83.2 | 62.5 | 66.0 | 85.1 | 80.7 | 84.9 | 27.2 | 73.2 | 57.5 | 78.1 | 79.2 | 81.1 | 77.1 | 53.6 | 74.0 | 49.2 | 71.7 | 63.3 |
| EdgeNet [58] | 71.2 | 83.6 | 35.8 | 82.4 | 63.1 | 68.9 | 86.2 | 79.6 | 84.7 | 31.8 | 74.2 | 61.1 | 79.6 | 76.6 | 83.2 | 80.9 | 58.3 | 82.6 | 49.1 | 74.8 | 65.1 |
| Attention [44] | 71.5 | 86.0 | 38.8 | 78.2 | 63.1 | 70.2 | 89.6 | 84.1 | 82.9 | 29.4 | 75.2 | 58.7 | 79.3 | 78.4 | 83.9 | 80.3 | 53.5 | 82.6 | 51.5 | 79.2 | 64.2 |
| DeepLab-Large [45] | 71.6 | 84.4 | **54.5** | 81.5 | 63.6 | 65.9 | 85.1 | 79.1 | 83.4 | 30.7 | 74.1 | 59.8 | 79.0 | 76.1 | 83.2 | 80.8 | 59.7 | 82.2 | 50.4 | 73.1 | 63.7 |
| CRFasRNN [59] | 72.0 | 87.5 | 39.0 | 79.7 | 64.2 | 68.3 | 87.6 | 80.8 | 84.4 | 30.4 | 78.2 | 60.4 | 80.5 | 77.8 | 83.1 | 80.6 | 59.5 | 82.8 | 47.8 | 78.3 | 67.1 |
| DeconvNet [60] | 72.5 | 89.9 | 39.3 | 79.7 | 63.9 | 68.2 | 87.4 | 81.2 | 86.1 | 28.5 | 77.0 | 62.0 | 79.0 | 80.3 | 83.6 | 80.2 | 58.8 | 83.4 | 54.3 | 80.7 | 65.0 |
| DPN [61] | 74.1 | 87.7 | 59.4 | 78.4 | 64.9 | 70.3 | 89.3 | 83.5 | 86.1 | 31.7 | 79.9 | 62.6 | 81.9 | 80.0 | 83.5 | 82.3 | 60.5 | 83.2 | 53.4 | 77.9 | 65.0 |
| Cont-CNN-CRF [46] | 75.3 | 90.6 | 37.6 | 80.0 | **67.8** | **74.4** | 92.0 | 85.2 | 86.2 | **39.1** | 81.2 | 58.9 | 83.8 | 83.9 | 84.3 | **84.8** | 62.1 | 83.2 | **58.2** | 80.8 | **72.3** |
| MoE-SPNet | 74.7 | 90.1 | 38.6 | 79.7 | 63.4 | 69.9 | 90.9 | 86.4 | **89.1** | 32.2 | **82.7** | 62.6 | **84.9** | 83.3 | **85.7** | 82.7 | **63.9** | 84.2 | 56.6 | 79.3 | 67.6 |
| FCN-AHFA | 70.6 | 82.6 | 37.2 | 80.9 | 58.0 | 67.7 | 86.4 | 84.6 | 84.5 | 30.2 | 76.6 | 50.3 | 78.7 | 79.1 | 83.4 | 80.3 | 59.3 | 78.5 | 48.5 | 80.5 | 61.9 |
| VGG + PASCAL VOC + COCO | | | | | | | | | | | | | | | | | | | | | |
| EdgeNet [58] | 73.6 | 88.3 | 37.0 | 89.8 | 63.6 | 70.3 | 87.3 | 82.0 | 87.6 | 31.1 | 79.0 | 61.9 | 81.6 | 80.4 | 84.5 | 83.3 | 58.4 | 86.1 | 55.9 | 78.2 | 65.4 |
| CRFasRNN [59] | 74.7 | 90.4 | 55.3 | 88.7 | 68.4 | 69.8 | 88.3 | 82.4 | 85.1 | 32.6 | 78.5 | 64.4 | 79.6 | 81.9 | 86.4 | 81.8 | 58.6 | 82.4 | 53.5 | 77.4 | 70.1 |
| BoxSup [47] | 75.2 | 89.8 | 38.0 | 89.2 | 68.9 | 68.0 | 89.6 | 83.0 | 87.7 | 34.4 | 83.6 | 67.1 | 81.5 | 83.7 | 85.2 | 83.5 | 58.6 | 84.9 | 55.8 | 81.2 | 70.7 |
| SBound [62] | 75.7 | 90.3 | 37.9 | 89.6 | 67.8 | 74.6 | 89.3 | 84.1 | 89.1 | 35.8 | 83.6 | 66.2 | 82.9 | 81.7 | 85.6 | 84.6 | 60.3 | 84.8 | 60.7 | 78.3 | 68.3 |
| Attention [44] | 76.3 | 93.2 | 41.7 | 88.0 | 61.7 | 74.9 | 92.9 | 84.5 | 90.4 | 33.0 | 82.8 | 63.2 | 84.5 | 85.0 | 87.2 | 85.7 | 60.5 | 87.7 | 57.8 | 84.3 | 68.2 |
| DPN [61] | 77.5 | 89.0 | **61.8** | 87.7 | 66.8 | 74.7 | 91.2 | 84.3 | 87.6 | 36.5 | 86.3 | 66.1 | 84.4 | 87.8 | 85.6 | 85.4 | 63.6 | 87.3 | 61.3 | 79.4 | 66.4 |
| Cont-CNN-CRF [46] | 77.8 | **94.1** | 40.4 | 83.6 | 67.3 | 75.6 | 93.4 | 84.4 | 88.7 | 41.6 | 86.4 | 63.3 | 85.5 | 89.3 | 85.6 | 86.0 | **67.4** | **90.1** | 62.6 | 80.9 | 72.5 |
| TVG-HO-CRF [63] | 77.9 | 92.5 | 59.1 | **90.3** | **70.6** | 74.4 | 92.4 | 84.1 | 88.3 | 36.8 | 85.6 | 67.1 | 85.1 | 86.9 | 88.2 | 82.6 | 62.6 | 85.0 | 56.3 | 81.9 | **72.5** |
| Att-CRF-DT [58] | 76.3 | 93.2 | 41.7 | 88.0 | 61.7 | 74.9 | 92.9 | 84.5 | 90.4 | 33.0 | 82.8 | 63.2 | 84.5 | 85.0 | 87.2 | 85.7 | 60.5 | 87.7 | 57.8 | 84.3 | 68.2 |
| MoE-SPNet | 77.7 | 91.6 | 39.7 | 89.6 | 64.2 | **77.1** | 93.7 | **89.0** | **93.6** | 36.5 | **87.6** | 56.0 | **90.3** | **91.6** | 85.9 | 86.7 | 59.2 | 89.3 | 59.3 | **85.7** | 70.9 |
| ResNet-101 + PASCAL VOC + COCO | | | | | | | | | | | | | | | | | | | | | |
| Deeplab-ASPP [8] | 79.7 | 92.6 | 60.4 | 91.6 | 63.4 | 76.3 | 95.0 | 88.4 | 92.6 | 32.7 | 88.5 | 67.6 | 89.6 | 92.1 | 87.0 | 87.4 | 63.3 | 88.3 | 60.0 | 86.8 | 74.5 |
| LRR-CRF [61] | 79.3 | 92.4 | 45.1 | 94.6 | 65.2 | 75.8 | 95.1 | 89.1 | 92.3 | 39.0 | 85.7 | 70.4 | 88.6 | 89.4 | 88.6 | 86.6 | 65.8 | 86.2 | 57.4 | 85.7 | 77.3 |
| Deep G-CRF [64] | 80.2 | 92.9 | 61.2 | 91.0 | 66.3 | 77.7 | **95.3** | 88.9 | 92.4 | 33.8 | 88.4 | 69.1 | 89.8 | 92.9 | 87.7 | 87.5 | 62.6 | 89.9 | 59.2 | **87.1** | 74.2 |
| FRRN [20] | 80.3 | 94.4 | 61.3 | 91.1 | 65.7 | 76.2 | 94.5 | 88.1 | 91.9 | 35.1 | 89.2 | 70.9 | 88.6 | 92.3 | 87.9 | 87.9 | 62.9 | 89.9 | 61.7 | 86.6 | 74.6 |
| Multi-Refine [65] | 82.4 | **94.9** | 60.2 | 92.8 | **77.5** | 81.5 | 95.0 | 87.4 | 93.3 | 39.6 | **89.3** | **73.0** | **92.7** | 92.4 | 85.4 | **88.3** | **69.7** | **92.2** | 65.3 | 84.2 | **78.7** |
| MoE-SPNet | **82.5** | 94.1 | **63.9** | **93.8** | 72.3 | **82.1** | 95.2 | **89.8** | **94.2** | **40.1** | 88.1 | 70.3 | 90.0 | **93.9** | **90.0** | 87.2 | 67.0 | 91.3 | **67.0** | 87.1 | 78.2 |

Table 1: **Results on PASCAL VOC 2012 test set.** For a fair comparison, In the bottom part of the table, we only compare results with previous works who also adopt the standard ResNet-101 as their base network. Thus, some works who modify the ResNet-101 to deeper or wider for their parsing network are not reported.
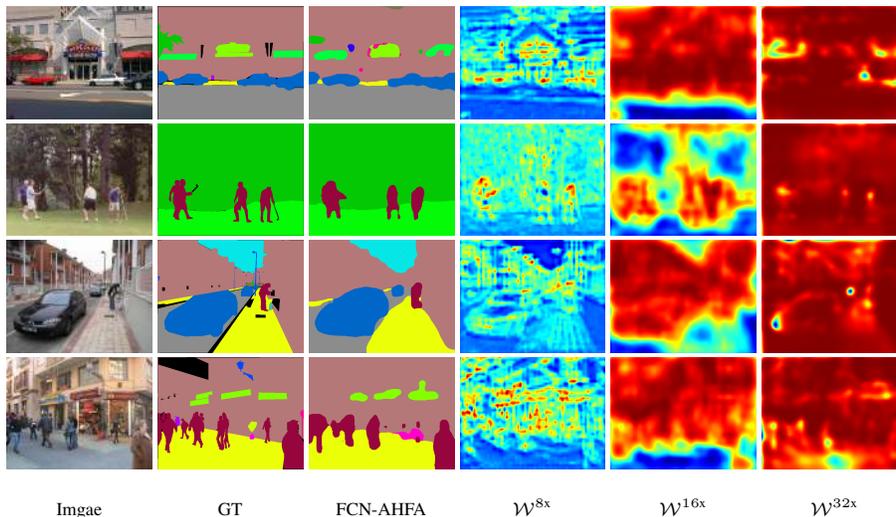
| Imgae | GT | FCN-AHFA | $\mathcal{W}^{8x}$ | $\mathcal{W}^{16x}$ | $\mathcal{W}^{32x}$ |

Figure 9: **Weight maps on SceneParse150 by FCN-AHFA.** Red represents high probability, and blue represents low probability. $\mathcal{W}^{8x}$, $\mathcal{W}^{16x}$, and $\mathcal{W}^{32x}$ correspond to weight maps of $\mathcal{F}^{8x}$, $\mathcal{F}^{16x}$, and $\mathcal{F}^{32x}$, respectively, as defined in section 4.2. $\mathcal{W}^{8x}$ give high weights to the boundary positions. $\mathcal{W}^{16x}$ sets high probabilities to regions of small categories for $\mathcal{F}^{16x}$. And the $\mathcal{W}^{32x}$ activates almost all the positions, since 32x features in FCN-8s capture the global contextual information which is useful for many categories. (Best viewed in colour)

### 5.2.2. SceneParse150

SceneParse150 [24] is a recently released large-scaled scene parsing benchmark using images from ADE20K Dataset [24]. The images are collected from diverse out-door and indoor scenes involving 35 stuff categories (*e.g.* floor, water, and sky) and 115 discrete objects (*e.g.* person, chair, and car). We train our models on the 20,210 training images, and evaluate the performance on the 2,000 validation images.

From the analyzation of different variants of MoE-SPNet, we choose MoE-SPNet, which contains least parameters but has promise performance, as our network architecture to demonstrate the effectiveness of our MoE based parsing algorithm. Also, following the same setting as the previous section, we employ the AHFA, which is the key component of MoE-SPNet, to the FCN. Note that, when building on ResNet, we add an extra convolutional layer with kernel size of $7 \times 7$ on top of ResNet to make

| Algorithm | Metric | | | |
|---|---|---|---|---|
| | Pixel Acc. | Mean Acc. | Mean IoU | Weighted IoU |
| VGG | | | | |
| Cascade-DilatedNet [24] | 74.52 % | 45.38 % | 0.3496 | 0.6108 |
| FCN-8s [5] | 71.56 % | 40.50 % | 0.2948 | 0.5755 |
| FCN-AHFA | **73.59** % | **43.51** % | **0.3128** | **0.6009** |
| DeepLab-ASPP [8] | 74.88 % | 46.17 % | 0.3303 | 0.6167 |
| MoE-SPNet | **75.50** % | **47.33** % | **0.3435** | **0.6242** |
| ResNet | | | | |
| FCN-16s [5] | 75.52 % | 44.13 % | 0.3475 | 0.6246 |
| FCN-AHFA | **76.04** % | **45.40** % | **0.3549** | **0.6286** |
| Deeplab-ASPP [8] | 77.31 % | 47.69 % | 0.3675 | 0.6354 |
| MoE-SPNet | **78.02** % | **48.02** % | **0.3789** | **0.6426** |

Table 2: **Results on SceneParse150 validation set.** We add AHFA and MOE to two kinds of parsing networks, including FCN-16/8s and Deeplab-ASPP, respectively. The comparison with baseline models demonstrates the effectiveness of our attention strategies.

the learned representation more complicated for FCN. We use FCN-16s as the baseline model, as we observed from the experiments that it outperforms FCN-8s on this dataset. Due to limited GPU memory, we use 50-layer ResNet for FCN-16s and FCN-AHFA and use ResNet101 to build two models based on DeepLab-ASPP and MoE-SPNet. Following the benchmark providers, we take the mean of *Pixel Acc.* and *Mean IoU* as the evaluation score.

The evaluation results are shown in Tab. 2. Sampled qualitative results are shown in Fig. 8. It can be seen that the AHFA-based methods outperform the baseline methods in terms of both pixel accuracy and IoU. Our VGG16-based FCN-AHFA yields a score of 52.4%, bringing 1.9% improvement over the VGG16-based FCN-8s (50.5%); and ResNet-based FCN-AHFA obtains a score of 55.8%, outperforming ResNet-based FCN-16s (55.1%) by 0.7%. Also, some selected learned weight maps in Fig. 9 furtherly demonstrate the effectiveness of our attention strategy.

Also, VGG-based MoE-SPNet has a score of 54.9%, yielding 0.9% improvement

over the baseline method DeepLab-ASPP. Also, ResNet-based MoE-SPNet achieves 58.0%, which is 1% higher than performance of ResNet-based DeepLab-ASPP. It should be noted that obtaining 1% overall improvement on this dataset containing 150 classes is considered as significant. Especially, the MoE-SPNet method outperforms the Cascade-DilatedNet [24] which segments stuff, objects, and object parts via a complicated cascade structure.

### 5.3. Ablation Studies on PASCAL VOC

We run some experiments to analyze our MOE and AHFA based networks, and discuss them in detail here.

### 5.3.1. Comparison with Baseline

| Algorithm | Mean IoU (%) | |
|---|---|---|
| | val | test |
| FCN-8s [5] | 68.4 | 62.2 |
| FCN-AHFA | 70.5 | 70.6 |
| DeepLab-ASPP [8] | 66.3 | 72.6 |
| MoE-SPNet | 70.4 | 74.7 |

Table 3: **Comparison with Baseline** This table reports the mean IoU on PASCAL VOC 2012 on val/test set. The FCN-8s and Deeplab-ASPP are two baseline networks.

To demonstrate the effectiveness of our MOE-SPNet, we compare VGG16-based MoE-SPNet with the baseline parsing network Deeplab-ASPP [8] on both the validation set and the test set. We also apply AHFA to the popular stage-wise parsing network FCN-8s to demonstrate the wide applicability of proposed AHFA based strategy. As shown in Tab. 3, our methods consistently outperform the counterpart baseline networks. In particular, our MoE-SPNet obtains 2.1% improvement in terms of *Mean IoU* compared with the baseline Deeplab-ASPP [8] on both sets. Significantly, employing our AHFA method to FCN-8s results in an *Mean IoU* of 70.4%, which not only outperforms FCN-8s (66.3%) by 4.1%, but also achieves comparable performance with

the state-of-the-art algorithms. It should be noted that, this comparison does not adopt the extra performance boosting techniques like CRF post-processing, pre-training on MS COCO, or multi-scale inputs.

### 5.3.2. Variants of MOE-SPNet

| | mean | areo | bike | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | mbike | person | plant | sheep | sofa | train | tv |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Baseline | 72.6 | 88.3 | 37.0 | 89.8 | 63.6 | 70.3 | 87.3 | 82.0 | 87.6 | 31.1 | 79.0 | 61.9 | 81.6 | 80.4 | 84.5 | 83.3 | 58.4 | 86.1 | 55.9 | 78.2 | 65.4 |
| MoE-SPNet-CF | 74.1 | 90.3 | 40.1 | 81.9 | 62.4 | 70.9 | 90.3 | 87.5 | 88.4 | 33.7 | 81.1 | 56.3 | 82.5 | 83.0 | 87.0 | 83.6 | 57.2 | 85.2 | 50.0 | 83.0 | 66.9 |
| MoE-SPNet-EF | 74.2 | 89.2 | 38.8 | 79.1 | 64.1 | 72.8 | 90.9 | 87.0 | 88.6 | 35.2 | 81.9 | 61.2 | 83.7 | 80.3 | 84.5 | 83.5 | 59.5 | 83.9 | 55.6 | 78.3 | 67.5 |
| MoE-SPNet | 74.7 | 90.1 | 38.6 | 79.7 | 63.4 | 69.9 | 90.9 | 86.4 | 89.1 | 32.2 | 82.7 | 62.6 | 84.9 | 83.3 | 85.7 | 82.7 | 63.9 | 84.2 | 56.6 | 79.3 | 67.6 |

Table 4: **Comparison of different MoEs for parsing.** Baseline: The DeepLab-ASPP parsing network without the gating part. MoE-SPNet-CF: The gating network using the input features to all the experts. MoE-SPNet-EF: The gating network takes the high-level features $\mathbb{F}$ within each expert as input. MoE-SPNet: The gating network takes the predictions $\mathcal{F}_i$ of each expert as input.

Furthermore, Tab. 4 shows evaluation results of variants of the proposed MoE-SPNet on the test server, including MoE-SPNet-CF whose gating network share the same input with the that of the experts, MoE-SPNet-EF whose gating network takes the features within each experts as input. From the table, all the MoE based parsing networks yield at least an improvement of 1.0% over the baseline network (Deeplab-ASPP). MoE-SPNet-EF outperforms MoE-SPNet-CF by 0.5%, while MoE-SPNet obtains a further 0.6% improvement compared with MoE-SPNet-EF, which demonstrate that direct understanding of a scene can help learn a more effective gating network. To exploit the MoE based parsing networks in deeper, we calculate the number of parameters of the gating networks belong to these MoE-SPNets here. Assuming the dimension of $\mathcal{S}$, $\mathbb{F}_i$, and $\mathcal{F}_i$ are $\mathcal{C}_1$, $\mathcal{C}_2$, and $\mathcal{C}_3$, respectively, and the additional convolutional layer for MoE-SPNet-EF and MoE-SPNet-CF contain $\mathcal{C}_4$ channels, the numbers of the gating networks parameters in MoE-SPNet-CF, MoE-SPNet-CF, MoE-SPNet are $\mathcal{C}_1 \times \mathcal{C}_4 \times 3 \times 3 + \mathcal{C}_4 \times \mathcal{N} \times 1 \times 1$, $\mathcal{N} \times \mathcal{C}_2 \times \mathcal{C}_4 \times 3 \times 3 + \mathcal{C}_4 \times \mathcal{N} \times 1 \times 1$, and $\mathcal{N} \times \mathcal{C}_3 \times \mathcal{N} \times 3 \times 3$, respectively, where $\mathcal{N}$ is the number of experts, and $\mathcal{N} < \mathcal{C}_3 \ll \mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_4$. It can be seen that MoE-SPNet contains the fewest parameters but achieves the best performance

by using a gating network learned from the predictions of each expert.

## 6. Conclusion

In this paper, we have proposed MoE-SPNet and FCN-AHFA to better exploit the diversities of contextual information in multi-level features and the spatial inhomogeneity of a scene in CNN-based models for scene parsing, by learning to assess the importance of features from different levels at each spatial location, instead of aggregating such features via concatenation or linear combination as commonly done in previous methods. The proposed MoE-SPNet achieves better performance by incorporating a mixture-of-experts layer to assess the importance of features from different layers. The AHFA scheme inspired by MoE-SPNet is applicable to a variety of scene parsing networks that use skip connections to fuse multi-level features from different stages. The value of the proposed methods have been demonstrated by the consistent and remarkable performance increase in a number of experiments on two challenging benchmarks (PASCAL VOC 2012 and SceneParse150). In the future, we will continue investigating more effective and efficient methods to jointly make use of multi-level convolutional features in CNN-based models for scene parsing and other challenging computer vision problems.

## References

## References

[1] M. Heber, M. Godec, M. Rüther, P. M. Roth, H. Bischof, Segmentation-based tracking by support fusion, Computer Vision and Image Understanding 117 (6) (2013) 573–586.

[2] Y. Fu, J. Cheng, Z. Li, H. Lu, Saliency cuts: An automatic approach to object segmentation, in: Pattern Recognition, 2008. ICPR 2008. 19th International Conference on, IEEE, 2008, pp. 1–4.

[3] B. Leibe, A. Leonardis, B. Schiele, Robust object detection with interleaved categorization and segmentation, International journal of computer vision 77 (1-3) (2008) 259–289.

[4] B. Hariharan, P. Arbeláez, R. Girshick, J. Malik, Simultaneous detection and segmentation, in: ECCV, Springer, 2014, pp. 297–312.

[5] J. Long, E. Shelhamer, T. Darrell, Fully convolutional networks for semantic segmentation, in: CVPR, 2015, pp. 3431–3440.

[6] J. Shotton, J. Winn, C. Rother, A. Criminisi, Textonboost for image understanding: Multi-class object recognition and segmentation by jointly modeling texture, layout, and context, International Journal of Computer Vision 81 (1) (2009) 2–23.

[7] P. Krähenbühl, V. Koltun, Efficient inference in fully connected CRFs with gaussian edge potentials, in: NIPS, 2011.

[8] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. L. Yuille, Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs, arXiv preprint arXiv:1606.00915.

[9] C. J. Holder, T. P. Breckon, X. Wei, From on-road to off: transfer learning within a deep convolutional neural network for segmentation and classification of off-road scenes, in: European Conference on Computer Vision, Springer, 2016, pp. 149–162.

[10] M. Kampffmeyer, A.-B. Salberg, R. Jenssen, Semantic segmentation of small objects and modeling of uncertainty in urban remote sensing images using deep convolutional neural networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2016, pp. 1–9.

[11] Z. Deng, S. Todorovic, L. Jan Latecki, Semantic segmentation of rgbd images with mutex constraints, in: Proceedings of the IEEE international conference on computer vision, 2015, pp. 1733–1741.

[12] S. R. Bulo, G. Neuhold, P. Kontschieder, Loss max-pooling for semantic image segmentation, CVPR.

[13] B. Hariharan, P. Arbeláez, R. Girshick, J. Malik, Hypercolumns for object segmentation and fine-grained localization, IEEE Transactions on Pattern Analysis and Machine Intelligence (2013) 1915–1929.

[14] S. Yin, Y. Qian, M. Gong, Unsupervised hierarchical image segmentation through fuzzy entropy maximization, Pattern Recognition 68 (2017) 245–259.

[15] Q. Zhou, B. Zheng, W. Zhu, L. J. Latecki, Multi-scale context for scene labeling via flexible segmentation graph, Pattern Recognition 59 (2016) 312–324.

[16] Z. Wang, L. Wei, L. Wang, Y. Gao, W. Chen, D. Shen, Hierarchical vertex regression based segmentation of head and neck ct images for radiotherapy planning 27 (2018) 923–937.

[17] G. Passino, I. Patras, E. Izquierdo, Pyramidal model for image semantic segmentation, in: Pattern Recognition (ICPR), 2010 20th International Conference on, IEEE, 2010, pp. 1554–1557.

[18] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: MICCAI, 2015.

[19] V. Badrinarayanan, A. Kendall, R. Cipolla, Segnet: A deep convolutional encoder-decoder architecture for image segmentation, IEEE TPAMI 39 (12) (2017) 2481–2495.

[20] T. Pohlen, A. Hermans, M. Mathias, B. Leibe, Full-resolution residual networks for semantic segmentation in street scenes, in: CVPR, 2017.

[21] G. Ghiasi, C. C. Fowlkes, Laplacian pyramid reconstruction and refinement for semantic segmentation, in: ECCV, Springer, 2016, pp. 519–534.

[22] M. I. Jordan, R. A. Jacobs, Hierarchical mixtures of experts and the em algorithm, Neural computation 6 (2) (1994) 181–214.

[23] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman, The pascal visual object classes challenge: A retrospective, International Journal of Computer Vision 111 (1) (2015) 98–136.

[24] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, A. Torralba, Semantic understanding of scenes through the ADE20K dataset, arXiv preprint arXiv:1608.05442.

[25] J. Shi, J. Malik, Normalized cuts and image segmentation, IEEE Transactions on pattern analysis and machine intelligence 22 (8) (2000) 888–905.

[26] P. F. Felzenszwalb, D. P. Huttenlocher, Efficient graph-based image segmentation, International journal of computer vision 59 (2) (2004) 167–181.

[27] B. Liu, H. Cheng, J. Huang, J. Tian, X. Tang, J. Liu, Fully automatic and segmentation-robust classification of breast tumors based on local texture analysis of ultrasound images, Pattern Recognition 43 (1) (2010) 280–298.

[28] H. Permuter, J. Francos, I. Jermyn, A study of gaussian mixture models of color and texture features for image classification and segmentation, Pattern Recognition 39 (4) (2006) 695–706.

[29] F. Bergamasco, A. Albarelli, A. Torsello, M. Favaro, P. Zanuttigh, Pairwise similarities for scene segmentation combining color and depth data, in: Pattern Recognition (ICPR), 2012 21st International Conference on, IEEE, 2012, pp. 3565–3568.

[30] M. Xian, Y. Zhang, H. Cheng, Fully automatic segmentation of breast ultrasound images based on breast characteristics in space and frequency domains, Pattern Recognition 48 (2) (2015) 485–497.

[31] M. Unger, M. Werlberger, T. Pock, H. Bischof, Joint motion estimation and segmentation of complex scenes with label costs and occlusion modeling, in: CVPR, 2012.

[32] E. Zemene, M. Pelillo, Interactive image segmentation using constrained dominant sets, in: European Conference on Computer Vision, Springer, 2016, pp. 278–294.

[33] T. Elguebaly, N. Bouguila, A nonparametric bayesian approach for enhanced pedestrian detection and foreground segmentation, in: Computer Vision and Pattern Recognition Workshops (CVPRW), 2011 IEEE Computer Society Conference on, IEEE, 2011, pp. 21–26.

[34] Z. Tu, X. Bai, Auto-context and its application to high-level vision tasks and 3d brain image segmentation, IEEE Transactions on Pattern Analysis and Machine Intelligence 32 (10) (2010) 1744–1757.

[35] J. Shotton, M. Johnson, R. Cipolla, Semantic texton forests for image categorization and segmentation, in: CVPR, IEEE, 2008, pp. 1–8.

[36] D. Ravì, M. Bober, G. M. Farinella, M. Guarnera, S. Battiato, Semantic segmentation of images exploiting dct based features and random forest, Pattern Recognition 52 (2016) 260–273.

[37] B. Fulkerson, A. Vedaldi, S. Soatto, Class segmentation and object localization with superpixel neighborhoods, in: Computer Vision, 2009 IEEE 12th International Conference on, IEEE, 2009, pp. 670–677.

[38] X. He, R. S. Zemel, M. A. Carreira-Perpiñán, Multiscale conditional random fields for image labeling, in: CVPR, 2004.

[39] S. P. Chatzis, D. I. Kosmopoulos, P. Doliotis, A conditional random field-based model for joint sequence segmentation and classification, Pattern recognition 46 (6) (2013) 1569–1578.

[40] A. Robles-Kelly, E. R. Hancock, A probabilistic spectral framework for grouping and segmentation, Pattern Recognition 37 (7) (2004) 1387–1405.

[41] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, ICLR.

[42] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: CVPR, 2016, pp. 770–778.

[43] C. Farabet, C. Couprie, L. Najman, Y. LeCun, Learning hierarchical features for scene labeling, IEEE Transactions on Pattern Analysis and Machine Intelligence 35 (8) (2013) 1915–1929.

[44] L.-C. Chen, Y. Yang, J. Wang, W. Xu, A. L. Yuille, Attention to scale: Scale-aware semantic image segmentation, in: CVPR, 2016.

[45] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. L. Yuille, Semantic image segmentation with deep convolutional nets and fully connected CRFs, in: ICLR, 2014.

[46] G. Lin, C. Shen, I. Reid, et al., Efficient piecewise training of deep structured models for semantic segmentation, in: CVPR, 2016, pp. 3640–3649.

[47] J. Dai, K. He, J. Sun, Boxsup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation, in: ICCV, 2015, pp. 1635–1643.

[48] G. Papandreou, I. Kokkinos, P.-A. Savalle, Untangling local and global deformations in deep convolutional networks for image classification and sliding window detection, in: ICCV, 2015.

[49] D. Bahdanau, K. Cho, Y. Bengio, Neural machine translation by jointly learning to align and translate, in: ICLR, 2015.

[50] P. H. O. Pinheiro, R. Collobert, Recurrent convolutional neural networks for scene labeling, in: ICML, 2014, pp. 82–90.

[51] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, G. E. Hinton, Adaptive mixtures of local experts, Neural computation 3 (1) (1991) 79–87.

[52] M. Holschneider, R. Kronland-Martinet, J. Morlet, P. Tchamitchian, A real-time algorithm for signal analysis with the help of the wavelet transform, in: Wavelets, Springer, 1990, pp. 286–297.

[53] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, ImageNet: A Large-Scale Hierarchical Image Database, in: CVPR, 2009, pp. 248–255.

[54] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, T. Darrell, Caffe: Convolutional architecture for fast feature embedding, arXiv preprint arXiv:1408.5093.

[55] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, A. Zisserman, The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results, http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html.

[56] B. Hariharan, P. Arbeláez, L. Bourdev, S. Maji, J. Malik, Semantic contours from inverse detectors, in: ICCV, 2011, pp. 991–998.

[57] M. Mostajabi, P. Yadollahpour, G. Shakhnarovich, Feedforward semantic segmentation with zoom-out features, in: CVPR, 2015, pp. 3376–3385.

[58] L.-C. Chen, J. T. Barron, G. Papandreou, K. Murphy, A. L. Yuille, Semantic image segmentation with task-specific edge detection using cnns and a discriminatively trained domain transform, in: CVPR, 2016, pp. 4545–4554.

[59] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, P. H. Torr, Conditional random fields as recurrent neural networks, in: ICCV, 2015, pp. 1529–1537.

[60] H. Noh, S. Hong, B. Han, Learning deconvolution network for semantic segmentation, in: ICCV, 2015, pp. 1520–1528.

[61] Z. Liu, X. Li, P. Luo, C.-C. Loy, X. Tang, Semantic image segmentation via deep parsing network, in: ICCV, 2015, pp. 1377–1385.

[62] I. Kokkinos, Pushing the boundaries of boundary detection using deep learning, in: ICLR, 2015.

[63] A. Arnab, S. Jayasumana, S. Zheng, P. H. Torr, Higher order conditional random fields in deep neural networks, in: ECCV, Springer, 2016, pp. 524–540.

[64] S. Chandra, I. Kokkinos, Fast, exact and multi-scale inference for semantic image segmentation with deep gaussian crfs, in: ECCV, Springer, 2016, pp. 402–418.

[65] G. Lin, A. Milan, C. Shen, I. Reid, RefineNet: Multi-path refinement networks for high-resolution semantic segmentation, in: CVPR, 2017.